

Calibrating Environmental Engineering Models

David Ruppert

Cornell University

September 12, 2007

Project Team

Calibrating
Environmental
Engineering
Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- **Christine Shoemaker**, co-PI, Professor of Civil and Environmental Engineering
 - PhD in applied math
 - works in applied optimization and environmental engineering
- David Ruppert, co-PI
- **Nikolai Blizniouk**, PhD student in Operations Research
- other students and post-docs
 - Rommel Regis
 - Stefan Wild
 - Pradeep Mugunthan
 - Dillon Cowan
 - Yingxing Li

Why is Calibration Difficult?

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- Likelihood may be multimodal
- Non-Gaussian data
- Spatial and temporal correlations
- non-constant variance: some data are much less accurate than others
- Model is computationally expensive
 - May take minutes or even hours to evaluate the model for one set of parameter values

Our Approach

Calibrating Environmental Engineering Models

David Rupert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- uses optimization and radial basis function meta-model to speed computations
- fully Bayesian
- takes into account all parameter uncertainty
- “noise” model includes possible
 - correlation
 - non-Gaussian distribution
 - non-constant variance

Bayesian versus Frequentist Statistics

Calibrating
Environmental
Engineering
Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating model

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- I do have sympathy with the Bayesian philosophy, but
 - I use Bayesian methods mainly as a powerful tool for finding estimators with good frequentist properties
- In general, the effect of the prior is $O(n^{-1})$
 - Estimation error is $O_P(n^{-1/2})$
- In complex nonlinear problems, exact confidence intervals are not impossible
 - Monte Carlo studies typically show the posterior credible intervals are approximate confidence intervals

Advantages of MCMC

Calibrating
Environmental
Engineering
Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- Non-Bayesian methods often use
 - the central limit theorem
 - linearization
- These approximations can create errors that are larger than the effect of the prior in a Bayesian analysis
- Even the bootstrap is justified by asymptotics:
 - the empirical CDF converges to the true CDF

Deterministic component of model

- i th observation is

$$Y_i = (Y_{i,1}, \dots, Y_{i,d})^T$$

- in absence of noise:

$$Y_{i,j} = f_j(X_i, \beta)$$

- $f_j(\cdot)$ comes from scientific theory
- X_i is a covariate vector
- β contains the parameters of interest
- noise is modeled empirically

What noise characteristic can we expect?

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- spatial and temporal correlations
- non-Gaussian distributions: most measured quantities are non-negative
- non-constant variance: **variance usually depends on the mean**
 - elephants vary more than mice
 - mice vary more than fleas

Components of the noise model

Calibrating
Environmental
Engineering
Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

We modeled the noise via:

- **data transformation:** to model
 - non-Gaussian variation
 - non-constant noise variance
- **spatial-temporal correlation model**

Transform-both-sides model

- The **transform-both-sides** model is

$$h \{ Y_{i,j}, \lambda_j \} = h \{ f_j(X_i, \beta), \lambda_j \} + \epsilon_{i,j},$$

- equivalently

$$Y_{i,j} = h^{-1} [h \{ f_j(X_i, \beta), \lambda_j \} + \epsilon_{i,j}, \lambda_j]$$

- transforms both sides of the equation giving deterministic model
- **preserves the theoretical model**
- $\{h(\cdot, \lambda) : \lambda \in \Lambda\}$ is some transformation family

Transform-both-sides examples

- the **identity transformation** gives the usual nonlinear regression model
 - **additive Gaussian errors**
- if we use the **log transformation** then

$$Y_{i,j} = \exp[\log\{f_j(X_i, \beta)\} + \epsilon_{i,j}] = f_j(X_i, \beta) \exp(\epsilon_{i,j})$$

- **multiplicative, lognormal errors**
- if we use the **square root transformation**

$$Y_{i,j} = \left[\sqrt{f_j(X_i, \beta)} + \epsilon_{i,j} \right]^2$$

- **notice a problem?**

The Box-Cox family

- the most common transformation family is due to Box and Cox (1964):

$$\begin{aligned}h(y, \lambda) &= \frac{y^\lambda - 1}{\lambda} \text{ if } \lambda \neq 0 \\ &= \log(y) \text{ if } \lambda = 0\end{aligned}$$

The Box-Cox family

- the most common transformation family is due to Box and Cox (1964):

$$\begin{aligned}h(y, \lambda) &= \frac{y^\lambda - 1}{\lambda} \text{ if } \lambda \neq 0 \\ &= \log(y) \text{ if } \lambda = 0\end{aligned}$$

- technical problem:
 - does not map $(0, \infty)$ onto $(-\infty, \infty)$, except for $\lambda = 0$
 - so transformed response has a truncated normal distribution
 - this makes Bayesian inference more complex

COIL transformation family

- **CONvex combination of Identity and Log (COIL) family:**

$$h_C(y, \lambda) = \lambda y + (1 - \lambda) \log(y), \quad 0 \leq \lambda \leq 1.$$

- We restrict λ to $[0, 1)$, since $h_C(\cdot, 1)$ does not map $(0, \infty)$ to $(-\infty, \infty)$
- COIL can approximate Box-Cox:
 - For each $\lambda \in [0, 1)$ there are constants $\lambda' \in [0, 1)$ and $a, b \in \mathbb{R}$ such that

$$h_{BC}(y, \lambda) \approx a + b h_C(y, \lambda')$$

for a wide range of y values (verified empirically)

- The inverse $h_C^{-1}(\cdot, \lambda)$ does not have a closed form
 - evaluate by interpolation (fast)

Multivariate transformations

Calibrating
Environmental
Engineering
Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- Define

$$\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_d)^T$$

- and

$$h(y, \boldsymbol{\lambda}) = \{h(y_1, \lambda_1), \dots, h(y_d, \lambda_d)\}^T$$

Separable correlation model

- Define the noise vectors:
 - $\epsilon_i = (\epsilon_{i,1}, \dots, \epsilon_{i,d})^T = h\{Y_i, \boldsymbol{\lambda}\} - h\{f(X_i, \boldsymbol{\beta}), \boldsymbol{\lambda}\}$
 - $\epsilon_{\bullet,j} = (\epsilon_{1,j}, \dots, \epsilon_{n,j})^T$
 - $\boldsymbol{\epsilon} = (\epsilon_1^T, \dots, \epsilon_n^T)^T$
- $\text{cov}(\epsilon_{i,j}, \epsilon_{i',j'}) = \boldsymbol{C}_{j,j'} \cdot \rho_{ST}(X_i, X_{i'}; \boldsymbol{\gamma})$
 - \boldsymbol{C} is a $d \times d$ covariance matrix for ϵ_i
 - $\rho_{ST}(X_i, X_{i'}; \boldsymbol{\gamma})$ is a space-time correlation function parameterized by $\boldsymbol{\gamma}$
- $\text{Var}\{\boldsymbol{\epsilon}\} = \boldsymbol{\Sigma}(\boldsymbol{\theta}) = \boldsymbol{S}(\boldsymbol{\gamma}) \otimes \boldsymbol{C}$
 - $\boldsymbol{\theta} = (\boldsymbol{\gamma}, \boldsymbol{C})$
 - $\boldsymbol{S}_{i,i'}(\boldsymbol{\gamma}) = \rho_{ST}(X_i, X_{i'}; \boldsymbol{\gamma})$

TBS Likelihood

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- Our statistical model is
$$h\{\mathbf{Y}, \boldsymbol{\lambda}\} \sim MVN[h\{\mathbf{f}(\boldsymbol{\beta}), \boldsymbol{\lambda}\}, \boldsymbol{\Sigma}(\boldsymbol{\theta})]$$
- Likelihood is

$$\begin{aligned} & [\mathbf{Y}|\boldsymbol{\beta}, \boldsymbol{\lambda}, \boldsymbol{\theta}] = \\ & \frac{\exp\left[-0.5 \|\mathbf{h}(\mathbf{Y}, \boldsymbol{\lambda}) - \mathbf{h}\{\mathbf{f}(\boldsymbol{\beta}), \boldsymbol{\lambda}\}\|_{\boldsymbol{\Sigma}(\boldsymbol{\theta})^{-1}}^2\right]}{(2\pi)^{nd/2} |\boldsymbol{\Sigma}(\boldsymbol{\theta})|^{1/2}} \cdot |J_h(\mathbf{Y}, \boldsymbol{\lambda})| \end{aligned}$$

- $|J_h(\mathbf{Y}, \boldsymbol{\lambda})|$ is the Jacobian
- $\boldsymbol{\Sigma}(\boldsymbol{\theta})$ is the covariance matrix

Overview of Methodology

- Goal:
 - Approximate the posterior density accurately with as few expensive likelihood evaluations as possible
- There are four steps:
 - 1 Locate the region(s) of high posterior density
 - 2 Find an “experimental design” that covers the region of high posterior density
 - the likelihood is evaluated on this design
 - 3 Use function evaluations from Steps 1 and 2 to approximate the posterior
 - 4 MCMC and standard Bayesian analysis using the **approximate** posterior density

Removing nuisance parameters

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- The posterior density is

$$[\beta, \lambda, \theta | \mathbf{Y}] = \frac{[\beta, \lambda, \theta, \mathbf{Y}]}{\int [\beta, \lambda, \theta, \mathbf{Y}] d\beta d\lambda d\theta},$$

- where $[\beta, \lambda, \theta, \mathbf{Y}] = [\mathbf{Y} | \beta, \lambda, \theta] \cdot [\beta, \lambda, \theta]$
- Interest focuses on

$$[\beta | \mathbf{Y}] = \int [\beta, \lambda, \theta | \mathbf{Y}] d\lambda d\theta$$

Removing nuisance parameters - four methods

- Exact:

$$[\beta | Y] = \int [\beta, \lambda, \theta | Y] d\lambda d\theta$$

- Profile posterior:

$$\pi_{\max}(\beta, Y) = \sup_{\zeta} [\beta, \zeta, Y] = [\beta, \hat{\zeta}(\beta), Y]$$

- $\hat{\zeta}(\beta)$ maximizes $[\beta, \zeta, Y]$ with respect to ζ
- Laplace approximation:
 - multiplies the profile posterior by a correction factor
- Pseudo-posterior:

$$[\beta, \hat{\zeta}(\hat{\beta}), Y]$$

- $\{\hat{\beta}, \hat{\zeta}(\hat{\beta})\}$ is the MAP = joint mode of posterior

Finding posterior mode using Condor

- When locating the posterior mode we want:
 - ① As few expensive function evaluations as possible
 - ② A small percentage of “wasted evaluations”
 - a) few evaluation locations in region of very low posterior probability
 - b) few evaluation locations that are very close together
 - ③ Getting very close to the mode is not a goal
- All good optimization techniques achieve 1
- Optimization methods based on numerical derivatives violate 2 b)
 - MATLAB's `fmincon` exhibited this problem
- CONDOR uses sequential quadratic programming
 - worked well in our empirical tests

Further function evaluations needed

- Goal:
 - approximate posterior on $C_R(\alpha) = \{\beta : [\beta, \mathbf{Y}] > \kappa(\alpha)\}$
- Function evaluations in optimization stage insufficient to approximate posterior accurately

Constructing the experimental design

① Normal approximation to posterior

- requires a small number of additional function evaluations

②

$$\hat{C}_R(\alpha) = \left\{ \beta : (\beta - \hat{\beta})^T [\hat{I}^{\beta\beta}]^{-1} (\beta - \hat{\beta}) \leq \chi_{p,1-\alpha}^2 \right\}$$

③ Space-filling design on $\hat{C}_R(\alpha)$

④ Remove points not in $\hat{C}_R(\alpha')$ for $\alpha' < \alpha$

- E.g., $\alpha = 0.1$ and $\alpha' = 0.01$

Radial basis functions

Calibrating
Environmental
Engineering
Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- $\pi(\cdot, \mathbf{Y})$ denotes one of the approximations to $[\beta, \mathbf{Y}]$

Radial basis functions

- $\pi(\cdot, \mathbf{Y})$ denotes one of the approximations to $[\boldsymbol{\beta}, \mathbf{Y}]$
- $l(\cdot) = \log\{\pi(\cdot, \mathbf{Y})\}$ is interpolated at $\mathcal{B}_D = \{\boldsymbol{\beta}^{(1)}, \dots, \boldsymbol{\beta}^{(N)}\}$ by

$$\tilde{l}(\boldsymbol{\beta}) = \sum_{i=1}^N a_i \phi(\|\boldsymbol{\beta} - \boldsymbol{\beta}^{(i)}\|_2) + q(\boldsymbol{\beta})$$

where

- $a_1, \dots, a_N \in \mathbb{R}$
- ϕ is a radial basis function
 - we used $\phi(r) = r^3$
- $q \in \Pi_m^p$ (the space of polynomials in \mathbb{R}^p of degree $\leq m$)
- $\boldsymbol{\beta} \in \mathbb{R}^p$

Autoregressive Metropolis-Hastings algorithm

- draw MCMC sample from $\tilde{\pi}(\cdot, \mathbf{Y}) = \exp\{\tilde{l}(\cdot)\}$
 - restrict sample to $\hat{C}_R(\alpha')$
- Metropolis-Hastings candidate:
$$\boldsymbol{\beta}^c = \boldsymbol{\mu} + \boldsymbol{\rho}(\boldsymbol{\beta}^{(t)} - \boldsymbol{\mu}) + \mathbf{e}_t$$
 - $\boldsymbol{\mu}$ = location parameter
 - $\boldsymbol{\rho}$ = autoregressive parameter (matrix)
 - $\rho = 0 \rightarrow$ independence MH
 - $\rho = 1 \rightarrow$ random-walk MH
 - \mathbf{e}_t 's are *i.i.d.* from density g
- if the candidate is accepted, then $\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^c$
- otherwise, $\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)}$

Applications in Environmental Engineering

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- few statisticians are working on environmental engineering problems
- environmental engineers typically use ad hoc and inefficient statistical methods
- modern statistical techniques such as variance functions, transformations, spatial-temporal models potentially offer substantial improvements
- statisticians and environmental will both benefit from collaboration

GLUE

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating model

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- GLUE = Generalized Likelihood Uncertainty Estimation
- widely used
- apparently considered state-of-the-art by many environmental engineers
- replaces the likelihood function of iid normal errors with an arbitrary objective function
- shows no appreciation of maximum likelihood as a general method
- objective function is not based on the data-generating probability model

Synthetic data example: Chemical spill

- To test algorithm:
 - use computationally inexpensive function
 - then approximate and exact result can be compared
- chemical accident caused spill at two locations on a long channel
 - mass M spill at location 0 at time 0
 - mass M spill at location L and time τ
- diffusion coefficient is d
- parameter vector is $\beta = (m, d, l, \tau)^T$
- want estimate of average concentration at end of channel
- l is of special interest
- need assessments of uncertainty as well

Chemical spill model

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- Model is:

$$\begin{aligned} C(s, t; M, D, L, \tau) &= \frac{M}{\sqrt{4\pi Dt}} \exp\left[\frac{-s^2}{4Dt}\right] \\ &+ \frac{M}{\sqrt{4\pi D(t-\tau)}} \exp\left[\frac{-(s-L)^2}{4D(t-\tau)}\right] \cdot \mathbb{I}(\tau < t) \end{aligned}$$

Details of simulation

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- assume data is collected at spatial location 0 (0.5) 2.5 and times 0.3 (0.3) 60 (5 time 200 observations)
- assume that a major goal is to estimate average concentration of time interval [40, 140] at the end of the channel ($s = 3$), specifically

$$F(\beta) = \sum_{i=0}^{20} f\{(3, 40 + 5i), \beta\}$$

- requires additional function evaluations (but not much more computation)

Details, continued

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- $\lambda = 0.333$ in COIL family
- one chemical species
- σ can be integrated out of the posterior analytically

Posterior densities: components of β

Calibrating
Environmental
Engineering
Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

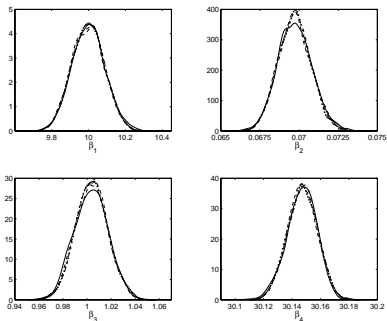


Figure: Kernel estimates of the posterior densities of β_i 's with the exact joint posterior (solid line) and RBF approximations to joint posterior (dashed line), pseudoposterior (dashed-dotted line), profile posterior with and without Laplace correction (dotted and large dotted lines, respectively).

Posterior densities: $F(\beta)$

Calibrating
Environmental
Engineering
Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

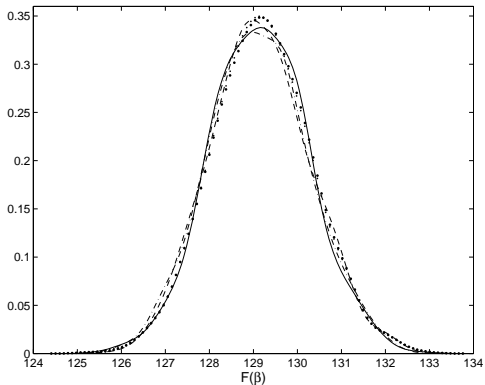


Figure: Kernel smoothed density estimates for the posterior of $F(\beta)$.

Results of a Monte Carlo experiment

Table: Observed coverage probabilities of Bayesian credible intervals.

	size .9 cred. int.		size .95 cred. int.		size .99 cred. int.	
	exact	RBF	exact	RBF	exact	RBF
β_1	.905 (.009)	.904 (.009)	.950 (.007)	.944 (.007)	.986 (.004)	.990 (.003)
β_2	.908 (.009)	.903 (.009)	.954 (.007)	.951 (.007)	.991 (.003)	.987 (.004)
β_3	.916 (.009)	.899 (.010)	.953 (.007)	.954 (.007)	.989 (.003)	.988 (.003)
β_4	.904 (.009)	.909 (.009)	.947 (.007)	.945 (.007)	.988 (.003)	.987 (.004)
$F(\beta)$.904 (.009)	.902 (.009)	.947 (.007)	.937 (.008)	.994 (.002)	.980 (.004)

What have we achieved?

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

In this research we have:

- applied modern statistical tools to calibration of environmental engineering models, e.g.,
 - transform-both-side
 - spatial-temporal correlation models
 - MCMC
- careful modeling of the noise increases estimation accuracy, often by a substantial amount
- implemented a Bayesian method of uncertainty analysis
- substantially reduced the number of evaluations of the computationally expensive environmental model by a meta-model based on RBF's

Current and Future Work

Calibrating Environmental Engineering Models

David Ruppert

Background

The team

The research problem

The Model

Environmental model

Modeling the noise

Likelihood

Methodology

Overview

Locating mode

Experimental Design

RBF approximation

MCMC sampling

Case Study

Chemical spill model

Monte Carlo

Summary

- multivariate observations, e.g., several chemical species
- multimodal posterior density
- covariate measurement error:
 - e. g., sampling error for rainfall can induce large correlated errors in a stream flow model
 - unlike response measurement error, covariate measurement error induces bias
- automatic tuning of MCMC